WORKSHOP 1

ANA CRISTINA QUINTERO CARPINTERO

2226286

Javier Alejandro Vergara Zorrilla

ETL
UNIVERSIDAD AUTÓNOMA DE OCCIDENTE

AGOSTO - 28 – 2024

WHAT IS EXPECTED?

I expect you to get the CSV file and create an application to migrate the data to a relational database. Also, you will display that data from the database in graphical visualizations - remember, the data should be stored in a database and your reports should come from the database, not the CSV file.

The visualizations I expect are:
Hires by technology (pie chart)
Hires by year (horizontal bar chart)
Hires by seniority (bar chart)
Hires by country over years (US, Brazil, Colombia, and Ecuador only) (multi-line chart)

Technologies
We expect you to use in this challenge:
Python
Jupiter Notebook
Database (Postgres)
Diagram

Workshop -001: Data Engineer 2
Data
I have 50,000 rows of data about candidates. The fields we'll be using are:
First Name
Last Name
Email
Country
Application Date
Year of Experience
Seniority
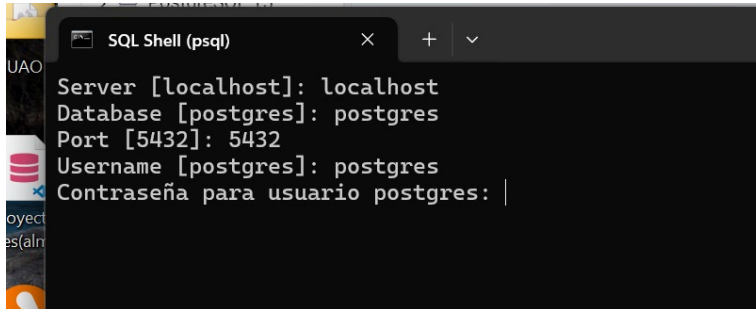Technology
Code Challenge Score
Technical Interview

Remember that I consider a candidate to be HIRED when they have both scores greater than or equal to 7 - you need to apply this logic to get the right information. How you'll handle this data is up to you.

And remember that all the data here is totally random - we used a public library to generate random information.

In this document, the entire development of workshop 1 will be presented.
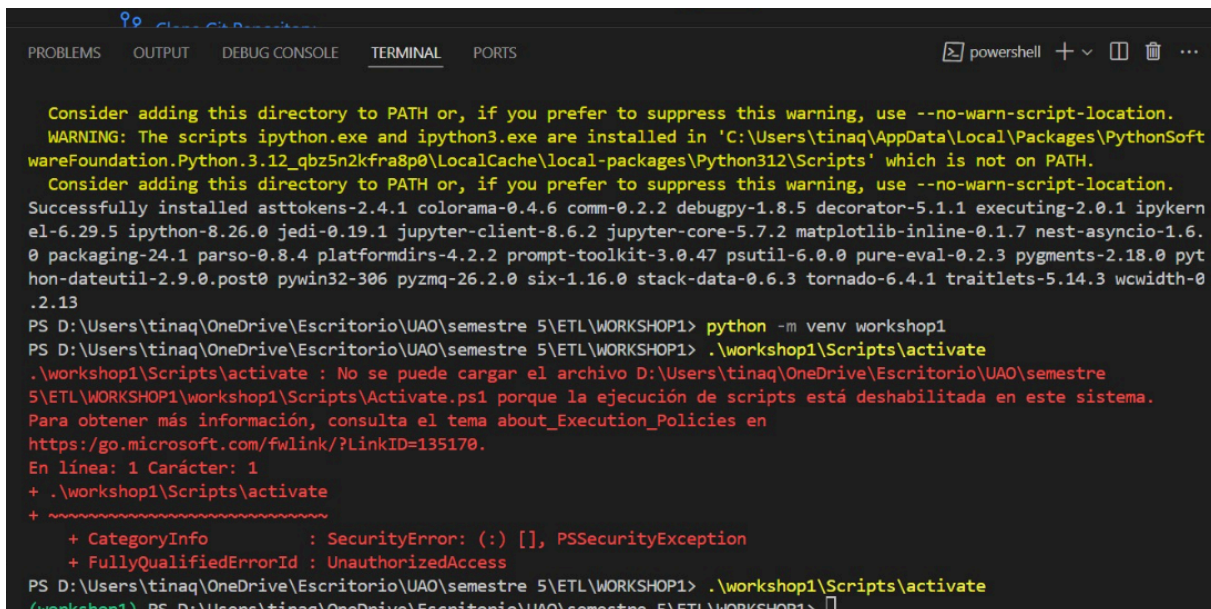
It is important to mention that the Postgres database was chosen, therefore, the connection to Postgres is created.



We proceed to create the database.



We proceed to create the kernel in the virtual environment



(workshop1) PS D:\Users\tinaq\OneDrive\Escritorio\UAO\semestre 5\ETL\WORKSHOP1>
pip install ipykernel

We proceed to install the libraries in the environment variable

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS    JUPYTER

● PS D:\Users\tinaq\OneDrive\Escritorio\UAO\semestre 5\ETL\WORKSHOP1> pip install python-dotenv pyt
  hon-decouple
  Collecting python-dotenv
    Downloading python_dotenv-1.0.1-py3-none-any.whl.metadata (23 kB)
  Collecting python-decouple
    Downloading python_decouple-3.8-py3-none-any.whl.metadata (14 kB)
  Downloading python_dotenv-1.0.1-py3-none-any.whl (19 kB)
  Downloading python_decouple-3.8-py3-none-any.whl (9.9 kB)
  Installing collected packages: python-decouple, python-dotenv
  Successfully installed python-decouple-3.8 python-dotenv-1.0.1
○ PS D:\Users\tinaq\OneDrive\Escritorio\UAO\semestre 5\ETL\WORKSHOP1>
```

We proceed to create the connection to the database

```
EDA.ipynb    .env    Extension: Python    pre_load.ipynb    db_conection.py ×    pyvenv.cfg

src > db_conection.py > ...
     1   from sqlalchemy import create_engine
     2   from decouple import config
     3
     4   engine = create_engine(f'postgresql://{config('DB_USER')}:{config('DB_PASSWORD')}@{config('DB_HOST'
     5
     6   class DbConnection:
     7       def _init_(self, eng=engine):
     8           self.engine = eng
     9
    10   conn = DbConnection()
```

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS    JUPYTER

Successfully installed python-decouple-3.8 python-dotenv-1.0.1
PS D:\Users\tinaq\OneDrive\Escritorio\UAO\semestre 5\ETL\WORKSHOP1> pip install sqlalchemy
Collecting sqlalchemy
  Downloading SQLAlchemy-2.0.32-cp312-cp312-win_amd64.whl.metadata (9.8 kB)
Collecting typing-extensions>=4.6.0 (from sqlalchemy)
  Downloading typing_extensions-4.12.2-py3-none-any.whl.metadata (3.0 kB)
Collecting greenlet!=0.4.17 (from sqlalchemy)
```

We proceed to verify that the database is previously loaded

```
SQL Shell (psql)          ×    + ∨

Esquema |      Nombre      |  Tipo  |  Dueño
--------+------------------+--------+----------
 public | candidates_raw   | tabla  | postgres
(1 fila)


workshop_1=# \d candidates_raw
                          Tabla «public.candidates_raw»
                Columna                 | Tipo | Ordenamiento | Nulable | Por omisión
----------------------------------------+------+--------------+---------+-------------
 First Name;Last Name;Email;Application Date;Country;YOE;Seniori | text |              |         |


workshop_1=# \d candidates_raw
                    Tabla «public.candidates_raw»
        Columna          |  Tipo  | Ordenamiento | Nulable | Por omisión
-------------------------+--------+--------------+---------+-------------
 First Name              | text   |              |         |
 Last Name               | text   |              |         |
 Email                   | text   |              |         |
 Application Date        | text   |              |         |
 Country                 | text   |              |         |
 YOE                     | bigint |              |         |
 Seniority               | text   |              |         |
 Technology              | text   |              |         |
 Code Challenge Score    | bigint |              |         |
 Technical Interview Score | bigint |            |         |
```
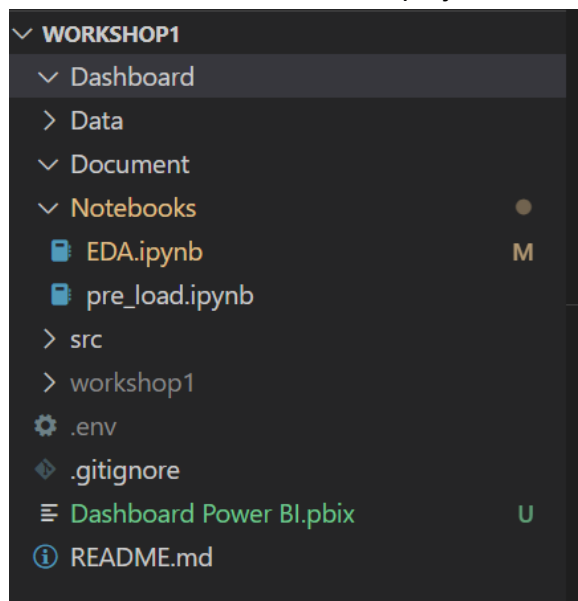
We validate that the dirty database is previously loaded.



We visualize the folders in the project:

Then we start the analysis process:
It is important to install the libraries and read the database

| First Name | Last Name | Email | Application Date | Country | YOE | Seniority | Technology | Code Challenge Score | Technical Interview Score |
|---|---|---|---|---|---|---|---|---|---|
| Bernadette | Langworth | leonard91@yahoo.com | 2021-02-26 | Norway | 2 | Intern | Data Engineer | 3 | 3 |
| Camryn | Reynolds | zelda56@hotmail.com | 2021-09-09 | Panama | 10 | Intern | Data Engineer | 2 | 10 |

A review is made of how the data is found:

```
# Verificar si hay filas duplicadas
print(f"Duplicated rows: {dataframe_raw.duplicated().sum()}")
0.1s
```

To do this, we review the distribution of the data and obtain the following results:

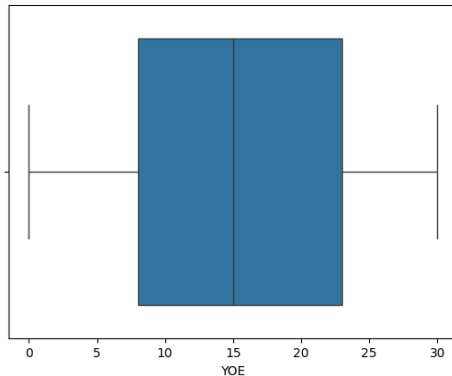| | YOE | Code Challenge Score | Technical Interview Score |
|---|---|---|---|
| count | 50000.000000 | 50000.000000 | 50000.000000 |
| mean | 15.286980 | 4.996400 | 5.003880 |
| std | 8.830652 | 3.166896 | 3.165082 |
| min | 0.000000 | 0.000000 | 0.000000 |
| 25% | 8.000000 | 2.000000 | 2.000000 |
| 50% | 15.000000 | 5.000000 | 5.000000 |

Subsequently, the exploratory data analysis process is carried out and therefore we obtain the following conclusions:
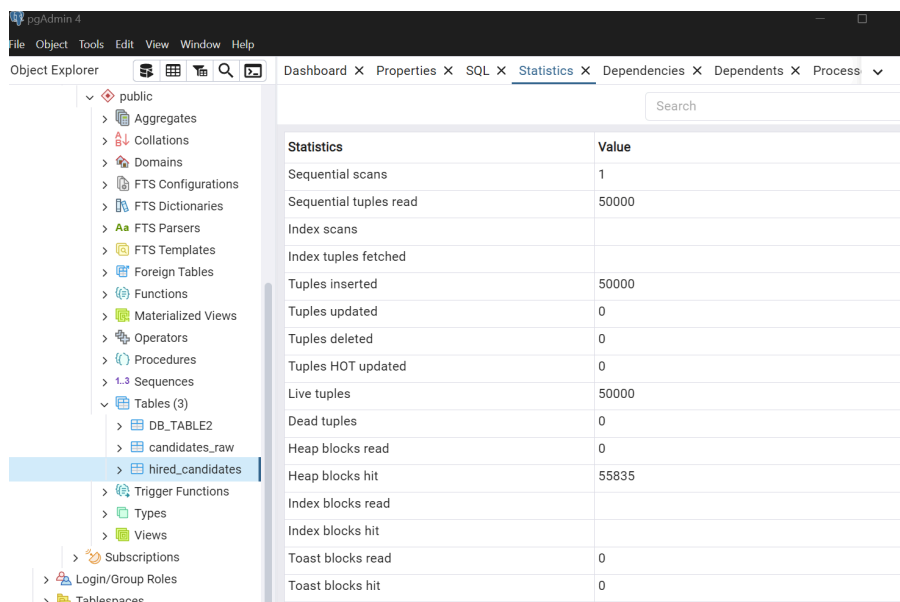


It is observed that there is a balanced distribution in the scores of the technical evaluations, with many candidates located in the middle range.

According to the graph, we can see that the years of experience variable shows a significant diversity in terms of the experience levels of the candidates.
It is detected that there is a large concentration of candidates with 30 years of experience and this indicates that there may be a concentration among the data, which leads us to look for outliers.



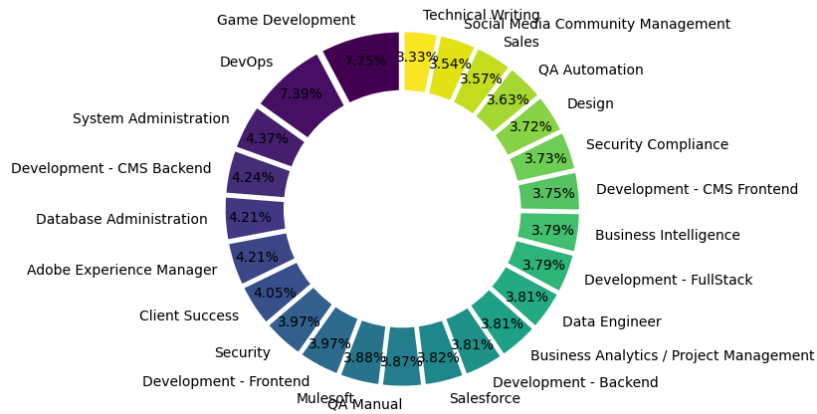After cleaning and EDA is saved in the database (hired_candidates)



Finally, since we have the clean table stored in our database, we begin to create the requested graphics.
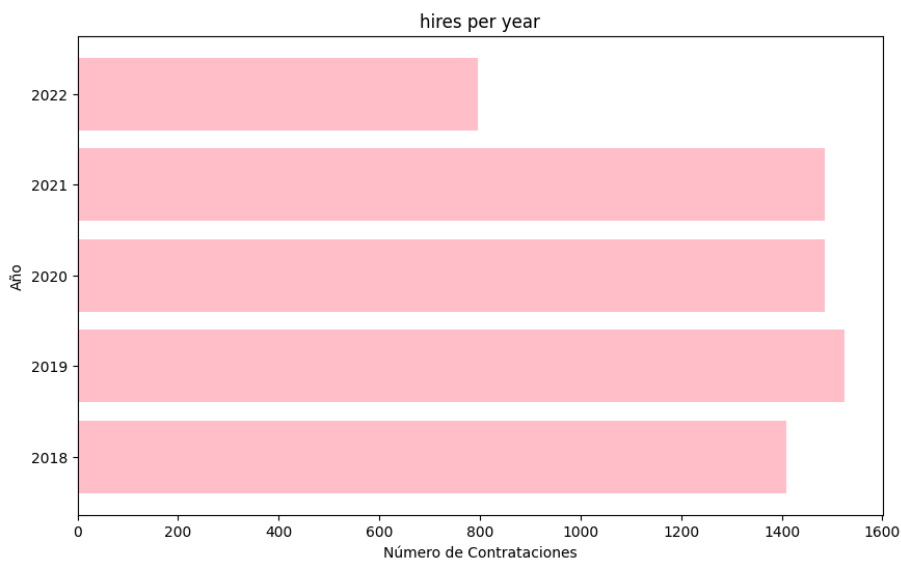
● Hiring by technology (pie chart)

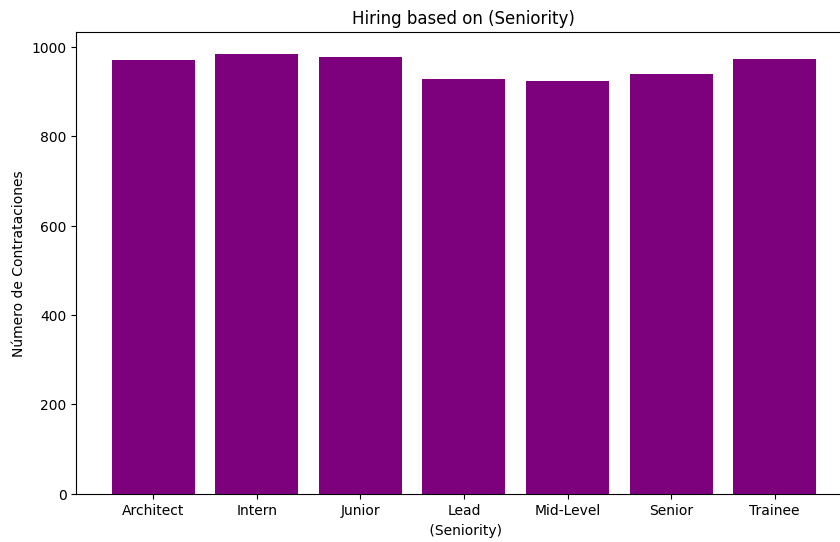**Distribución de Contrataciones por Tecnología**



In the graph above we can see how the hiring of prospective employees is distributed across various fields in the technology area.

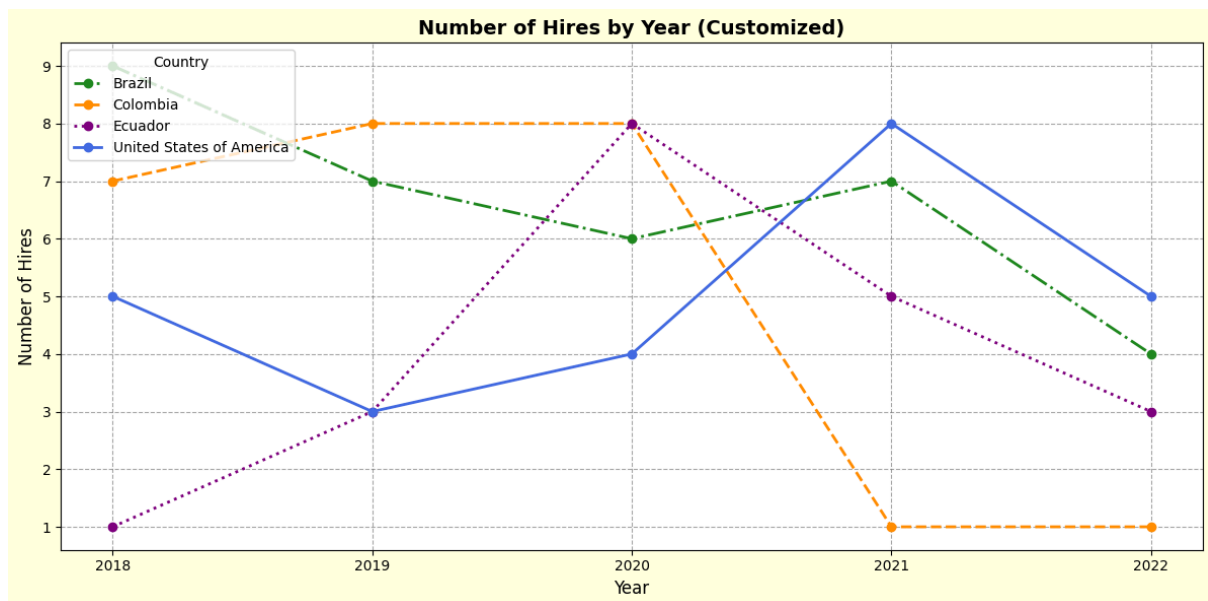● Hires per year (horizontal bar chart)

hires per year



In the graph above we can see how the hiring of prospects is distributed over the various years.

● Hires by seniority (bar chart)



In the graph above we can see how the hiring of prospective employees is distributed by seniority.

● Hires by country over the years (US, Brazil, Colombia and Ecuador only) (multi-line chart)

In the graph above we can see how the hiring of prospects is distributed by selected countries over the years.

Finally, the database connection to the Power BI tool is made.



The following business questions were asked to implement in Power BI.

Which countries have contributed the highest number of candidates?
How does the years of experience (YOE) vary based on the seniority level of the candidates?
Is there any correlation between the country of origin of the candidates and the technologies they are specialized in?
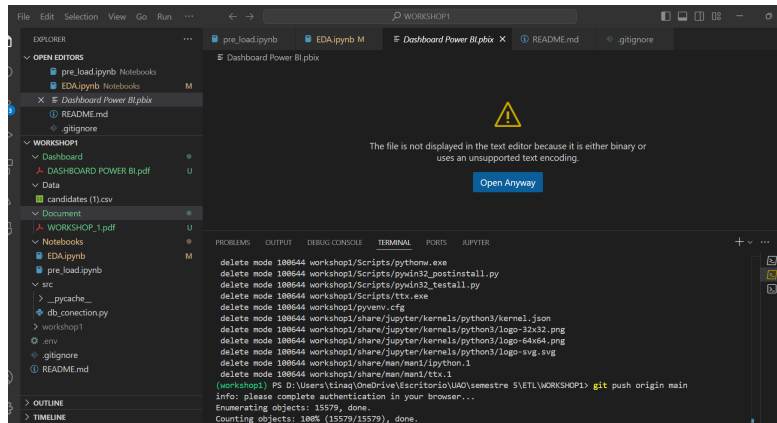What is the average duration between the application date and the hiring date?
Relationship between the year of application with the year of the technical interview
Relationship between Code Challenge Score and Technical Interview

Se sube desde el Visual Studio, todo el trabajo a GitHub



Los comandos usados fueron:

git init

**git rm --cached -r workshop1**

**git add .**

**git commit -m**

**git push origin main**